

(N)

PAT-NO: JP406215196A  
DOCUMENT-IDENTIFIER: JP 06215196 A  
TITLE: METHOD FOR ESTABLISHING GROUP DATA BASE IN CHARACTER  
RECOGNITION  
PUBN-DATE: August 5, 1994

INVENTOR-INFORMATION:  
NAME  
TO, RAKUTEI  
SEN, EIKAN  
JO, EIJI  
RIN, BUNBUN

ASSIGNEE-INFORMATION:  
NAME IND TECHNOL RES INST COUNTRY  
N/A

APPL-NO: JP04318780  
APPL-DATE: November 27, 1992

INT-CL (IPC): G06K009/68

ABSTRACT:

PURPOSE: To reduce overlapped amounts, and to improve group applying precision by applying a pattern clustering group center to a group.

CONSTITUTION: The plural leaning patterns of one character are ununiformly divided (ST1), and the characteristics of the learning patterns are extracted (ST2). Then, the constant number of features are selected from each feature of the learning patterns, and the capability of a group is obtained (ST3). The learning patterns are divided into  $\omega$  groups by a clustering method based on the selected features (ST4). This  $\omega$  is an integer larger than 1. Then, the center value of the  $\omega$  groups obtained by clustering the learning patterns is applied to the group as a central pattern, overlapping is obtained, and the best grouping is obtained (ST5). The clustering method includes a random clustering, ISODATA clustering, and K-mean clustering, and they are sequentially executed.

COPYRIGHT: (C)1994, JPO

(19)日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11)特許出願公開番号

特開平6-215196

(43)公開日 平成6年(1994)8月5日

(51)Int.Cl.<sup>5</sup>

G 0 6 K 9/68

識別記号

庁内整理番号

9289-5L

F I

技術表示箇所

審査請求 未請求 請求項の数3 O L (全 6 頁)

(21)出願番号 特願平4-318780

(22)出願日 平成4年(1992)11月27日

(71)出願人 390023582

財団法人工業技術研究院

台湾新竹縣竹東鎮中興路四段195號

(72)発明者 屠 樂 挺

台湾台北市士林區名山里5鄰▲雨▼聲街53巷2號4樓

(72)発明者 ▲せん▼ 永 寛

台湾彰化縣永靖鄉▲なん▼港村永社路260巷2鄰18號

(72)発明者 徐 英 士

台湾台北市北投區文林里16鄰致遠一路一段47巷1號2樓

(74)代理人 弁理士 伊東 忠彦 (外1名)

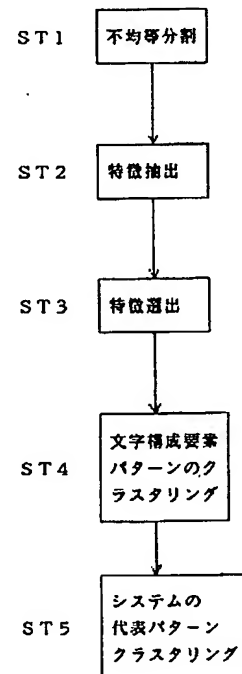
最終頁に続く

(54)【発明の名称】 文字認識におけるグループデータベース確立の方法

(57)【要約】

【目的】 文字認識において文字の代表パターンを利用してシステムとしてグループデータベースを確立する方法であって、パターンが多すぎること及び変異が大きすぎること起因するシステムのグループデータベースの拡張(重複量の増大)を回避すると同時に、システムのデータベースの保存スペースを削減する文字認識におけるグループデータベース確立方法を提供することを目的とする。

【構成】 1つの文字について多数個の学習パターンをそれぞれ不均等分割すること、学習パターンの特徴を抽出すること、学習パターンの各特徴の中から一定数の特徴を選出して、グループを確立すること、選出された特徴に基づいて、学習パターンをクラスタリングして $\omega$ グループ( $\omega$ は1より大きい整数)とすること、および、学習パターンを分けた $\omega$ グループの中心を代表パターンとして、クラスタリングによりグループ内部で重複を行って、最良のグループ分けを得ることの各ステップよりなる。



## 【特許請求の範囲】

【請求項1】 1つの文字の他数個の学習パターンをそれぞれ不均等分割すること、  
 当該文字の学習パターンの特徴を抽出すること、  
 各学習パターンの特徴から一定数の特徴を選出して、グループを構成すること、  
 選出した特徴に基づき、学習パターンをクラスタリング法で $\omega$ グループに分けるが、この $\omega$ が1より大きい整数であること、  
 および、  
 学習パターンをクラスタリングした $\omega$ グループの中心を代表パターンとして、クラスタリング法を利用してグループ内に投入して重複させ、最良のグループ分けを実現することという上記のステップを具備する文字認識におけるグループデータベース確立の方法。

【請求項2】 上記クラスタリング法が、ランダムクラスタリング、ISODATAクラスタリング、k-平均クラスタリングを含み、順序だてて実行されるものである請求項1記載のグループデータベース確立の方法。

【請求項3】 上記 $\omega$ が、 $\omega \geq 1.5$ である請求項1記載のグループデータベース確立の方法。

## 【発明の詳細な説明】

## 【0001】

【産業上の利用分野】本発明は文字認識におけるグループデータベース確立の方法に関し、特に文字認識において文字の代表パターンを利用してシステムとしてグループデータベースを確立する方法であって、パターンが多すぎることに及び変異が大きすぎることに起因するシステムのグループデータベースの拡張（重複量の増大）を回避すると同時に、システムのデータベースの保存スペースを削減する方法に係る。

## 【0002】

【従来の技術】従来において、文字認識の技術分野では、グループデータベース確立の技術について考慮を要する2点、つまりパターン重複量とグループ投入精度とがあった。一般に使用されてきた重複技術では、グループを先ず $C_1, C_2, \dots, C_m$ と仮設する。文字構成要素（訳注：字素？ 次元？） $X_i$ につき、 $1 \leq i \leq N$ 、 $X_i$ と $C_k$ との距離を $d(X_i, C_k) = \min d(X_i, C_k)$ としていた。もし、

## 【0003】

## 【数1】

$$X_i \in C_k$$

【0004】ならば $C_k \cup \{X_i\}$ が重複を示す。

## 【0005】

【発明が解決しようとする問題点】しかしながら、従来においては、重複量とグループ投入精度とは両立できないものであった。もしグループ投入精度を向上させたいなら重複量を増やさなければならぬので、グループ内

部に保存できるコードの保存スペースに影響をおよぼすとともに、マッチングに要する時間を増大させるものとなっていた。従って、このような欠点を克服することが、文字認識の技術分野における重要課題となっていた。

## 【0006】

【問題点を解決するための手段】上記した公知の問題を解決するために、本発明では、認識プロセスにおいて多数個の候補グループを利用することによって、従来、代表パターンが完全には字形の変異を吸収できなくなり、かえって1グループだけを探し出すために誤識別というエラーを引き起こしていたことを回避するとともに、累積により精度が増大するという現象を利用して、多グループを採用することによりグループ投入精度を向上させる。

【0007】このような発明の観点から、文字認識におけるグループデータベース確立の方法を提供するもので、下記ステップを具備する。すなわち、1つの文字について多数個の学習パターンをそれぞれ不均等分割すること、学習パターンの特徴を抽出すること、学習パターンの各特徴の中から一定数の特徴を選出して、グループを確立すること、選出された特徴に基づいて、学習パターンをクラスタリングして $\omega$ グループ（ $\omega$ は1より大きい整数）とすること、および学習パターンを分けた $\omega$ グループの中心を代表パターンとして、クラスタリングによりグループ内部に投入して重複を行って、最良のグループ分けを得ること、という各ステップである。

## 【0008】

【実施例】以下、図1のフローチャートを参照しながら、本発明に係るグループデータベース確立の方法を説明する。最初に、ステップST1において、文字の特徴値を抽出する前に、局部および内部字形のデータを獲得するために、先ず文字画像を分割する。文字の変形、例えば文字の高さや幅の不一致、字形や字画の不均等な分布、ならびに特殊な字画の長さの差異により、もし印刷文字方式の均等分割を行うならば、かならず分割が不正確となって、抽出した特徴値の変異が非常に大きいものとなり、誤識別を発生させる。そこで、本発明は、不均等分割方式を採用する。例をあげるなら、図4の『示』という文字を分割するために、3種類の不均等分割方式、つまり字画密度関数（Stroke Density Function; SDF）、周辺面積（Peripheral Background Area; PBA）、周辺輪郭線長さ（Contour Line Length; CLL）を採用し、それぞれ図5の（A）、図5の（B）、図5の（C）に示す。これら3種類の不均等分割方式は、いずれも公知技術なので、ここでは改めて説明しない。

【0009】そして、ステップST2において、本発明に係る実施例は、SDF、PBA、CLLという3種類

の不均衡分割で64個の特徴を取り出す。続いて、ステップST3において、

$$CF_j = \sum_k \frac{M_{ji} - m_{jk}}{\sigma_{jk}}$$

式中、jは特徴の第j次元、i, kは第i, kの文字、mは中心値、 $\sigma_{jk}$ は文字の第j次元の特徴標準偏差STDを表す、

【0011】を利用して、 $CF_j$  値 ( $j=1\sim64$ ) のなかで最大となる40個を探して、グループ化された40次元の特徴とするとともに、残りの24個を棄却して使用しない。次に、ステップST4において、文字構成要素のパターンクラスタリングを行う。図2を見ると分かりやすいように、本発明の好適な実施例に係るクラスタリング法(H-class)の流れを詳細に示している。本発明の好適な実施例において、2568個の文字構成要素を採用している。各文字構成要素に対して、多数個の(本実施例では100個の)学習パターンをクラスタリング法で先にωグループに分ける。ωは、本実施例では、15, 17, 20に相当する。クラスタリング法は、まずランダム・クラスタリングを行って、次にイソデータISODATAクラスタリングを1回、さらにK-mean(K平均)クラスタリングを行い、後者によって得られた中心値を文字構成要素を代表する代表パターンとする。そして、次の文字構成要素について10※

10※0個の学習パターンを分解して、2568個の文字構成要素の全部を処理するまで続ける。

【0012】上記イソデータクラスタリングは、 $X_1, X_2, \dots, X_n$  を  $C_1, C_2, \dots, C_m$  中に投入して、以下の4ステップを行う必要がある。

(1)  $X_i$  を任意に  $C_k$  グループ中へ分散して、 $count=0$  とする、(2) 各  $X_i$  につき、 $1 \leq i \leq N$ ,  $X_i \in C_s$  とし、 $d(X_i, C_t) = \min d(X_i, C_k)$ ,  $1 \leq k \leq m$  とする、(3) もし  $t \neq s$  なら、すなわち  $C_t + \{X_i\}$ ,  $C_s = C_s - \{X_i\}$ ,  $count = count + 1$  とする、(4)  $i = i + 1$  とし、もし  $i = N$  なら  $C_k$  中の成分について改めて中心値を計算し、違うなら(2)に戻る。

【0013】ここで、

【0014】

【数3】

$$d(X_i, C_k) = \frac{(X_{ij} - C_{kj})^2}{\sigma_{kj} \cdot \sigma_{ij}}$$

特徴値  $X_i = \begin{bmatrix} X_{i1} \\ \vdots \\ X_{iF} \end{bmatrix}$ 、中心値  $C_k = \begin{bmatrix} C_{k1} \\ \vdots \\ C_{kF} \end{bmatrix}$ 、Fは特徴の次元数、

$\sigma_k = \begin{bmatrix} \sigma_{k1} \\ \vdots \\ \sigma_{kF} \end{bmatrix}$  を第kグループの標準偏差STDとし、

$$\sigma_{ij} = \frac{1}{2568} \sum_{n=1}^{2568} \Delta_{nj} \text{ とするが、}$$

【0015】 $\Delta_{nj}$ は、第n字の第j特徴の標準偏差STD(注意：本実施例では、2568個の字を使用)とする。また、k-平均クラスタリングは、以下の5ステップを備えている。

(1)  $X_i$  を任意に  $C_k$  グループ中に分散して入れ、 $count=0$  とする、(2) 各  $X_i$  ごとに、 $1 \leq i \leq N$ ,  $X_i \in C_s$  とする、(3) もし  $t \neq s$  なら  $C_t = C_t + \{X_i\}$ ,  $C_s = C_s - \{X_i\}$ ,  $count = count + 1$  とする、(4)  $i = i + 1$  とし、 $C_k$  中の成分を利用して改めて中心値を計算する、(5) もし  $count < \text{臨界値} T$  なら停止し、そうでなければ  $count=0$  として、(2)に戻る。

【0016】2568個の文字構成要素の各学習パターンをすべてクラスタリングした後、ステップST5に進み、システムの代表パターンのクラスタリングを行う。ステップST5の詳細な流れは図3を参照されたい。こ

の時、各パターン毎に $k$ -平均クラスタリングでクラス分けした $\omega$ グループの中心値を代表パターンとしてグループ内に投入して、もう一度、上記のクラスタリング法を応用して重複させる。そして、パターンを測定して照合するとともに、数個の候補グループを選出する（本実施例では5個としている）。その後、各パターンごとにグループ投入精度が最大となる時の $\omega$ を取る。そして、もう一度、ステップST5を繰り返す。本実施例によれば\*

\*ば、2568字が200個のパターンを有するシステムにおいて、100個のパターンを学習パターンとし、他の100個は測定を行う。本発明により得られる識別効果は極めて顕著なものであるので、公知技術と比較した表1を参照されたい。

【0017】

【表1】

	公 知	本 発 明
字数	2 5 6 8	2 5 6 8
クラスタリング法	$k$ -平均	$k$ -平均
学習パターン数	100パターン/字	100パターン/字
重複方法	パターンをグループに投入	パターンクラスタリンググループ中心をグループに投入
重複量	10	2.83
候補グループ数	1	5
測定パターン数	100	100
各グループの字数平均	700	760
グループ投入精度	97.46%	99.58%

【0018】上の表から分かるように、本発明はグループ中心をグループに投入することによって、重複量を削減させるばかりではなく、グループ投入精度を向上させる。以上に記載した比較的好ましい実施例は、本発明を説明するためのもので、本発明を制限するものではない。この技術分野に詳しい当業者なら多くの修正とバリエーションとが可能であるが、それでもなお、以下の特許請求の範囲で限定した発明のカテゴリーと精神とを回避できるものではない。

【図面の簡単な説明】

【図1】本発明に係るグループデータベース確立の方法を示すフローチャートである。

【図2】図1のステップST4の詳細な流れを示すフロ※

※ーチャートである。

【図3】図1のステップST5の詳細な流れを示すフローチャートである。

【図4】1つの中国文字『示』の不均衡分割を示す図である。

【図5】図4の『示』という文字を3種類の不均衡分割方式で分割して、特徴抽出した図である。

【符号の説明】

ST1 不均等分割

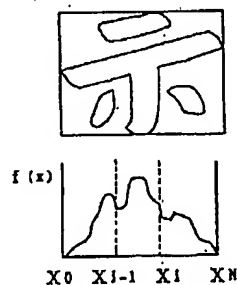
ST2 特徴抽出

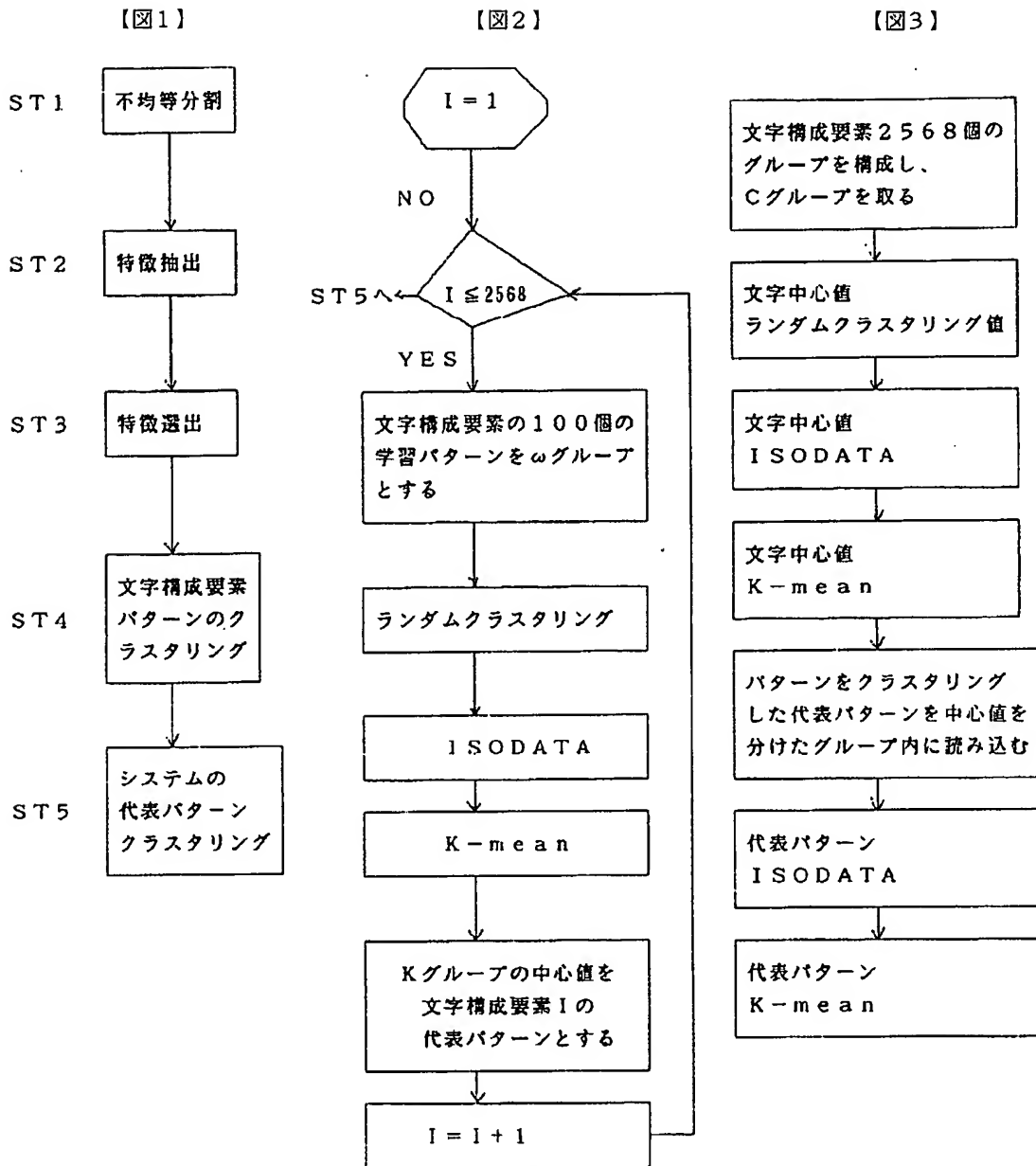
ST3 特徴選出

ST4 文字構成要素のクラスタリング

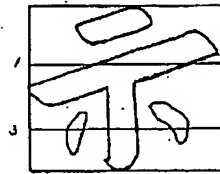
ST5 システムの代表パターンクラスタリング

【図4】

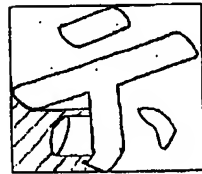




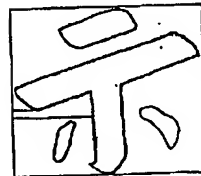
【図5】



(A)



(B)



(C)

---

フロントページの続き

(72)発明者 林 文 ▲ぶん▼  
台湾台北市北投區温泉里10鄰温泉路58巷9  
-1 號